

---

# EEVR: A Dataset of Paired Physiological Signals and Textual Descriptions for Joint Emotion Representation Learning

---

Pragya Singh<sup>1</sup>, Ritvik Budhiraja<sup>1</sup>, Ankush Gupta<sup>1</sup>, Anshul Goswami<sup>1</sup>, Mohan Kumar<sup>2</sup>, and Pushpendra Singh<sup>1</sup>

<sup>1</sup>IIIT-D, New Delhi, India, pragyas@iiitd.ac.in, ritvik19322@iiitd.ac.in, ankush21232@iiitd.ac.in, anshul20361@iiitd.ac.in, psingh@iiitd.ac.in

<sup>2</sup>RIT, Rochester, New York, US, mjkvcs@rit.edu

## 1 Supplementary Materials

2	1	EEVR Overview	3
3	1.1	EEVR Size Details . . . . .	3
4	1.2	EEVR Organization and File formats . . . . .	3
5	2	EEVR Usage and Publishing	5
6	2.1	Intended Uses . . . . .	5
7	2.2	Ethical Consideration . . . . .	5
8	2.3	EEVR Licensing, Hosting, and Maintenance Plan . . . . .	5
9	3	Human Subjects Considerations	6
10	4	Data Collection Protocol	6
11	4.1	Experiment Instruction and Sensors Preparation . . . . .	6
12	4.2	Stimulus Selection and Playlist Preparation . . . . .	6
13	4.3	Virtual Reality Module Preparation . . . . .	7
14	4.4	Self-Assessment . . . . .	8
15	4.4.1	GHQ-12 . . . . .	9
16	4.4.2	Personality . . . . .	9
17	4.4.3	VRSQ . . . . .	10
18	5	Data Analysis and Experiments	11
19	5.1	Content Analysis . . . . .	11
20	5.2	Data Cleaning . . . . .	11
21	5.3	Text Data Preparation . . . . .	12
22	5.4	Text Data Analysis . . . . .	12
23	5.5	Physiological Features . . . . .	12
24	5.6	Experiment Details and Results Analysis . . . . .	13
25	5.6.1	Experimental Setup . . . . .	13
26	5.6.2	Discussion on Physiological Baseline . . . . .	13
27	5.6.3	Discussion on Text Baseline . . . . .	13

Submitted to the 38th Conference on Neural Information Processing Systems (NeurIPS 2024) Track on Datasets and Benchmarks. Do not distribute.

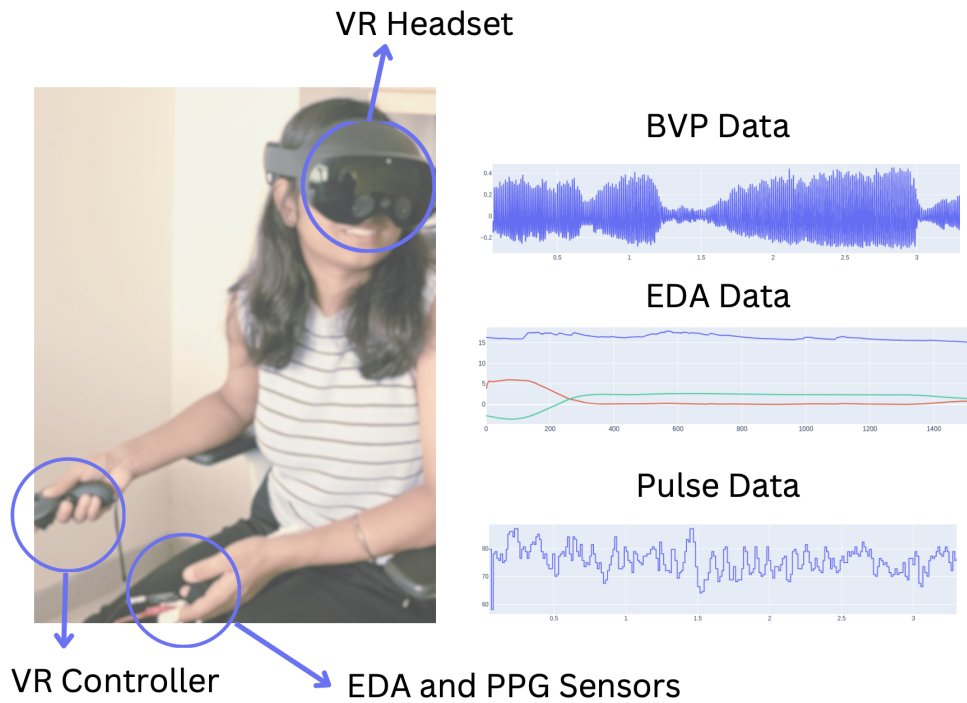


Figure 1: Experiment Setup of EEVR dataset

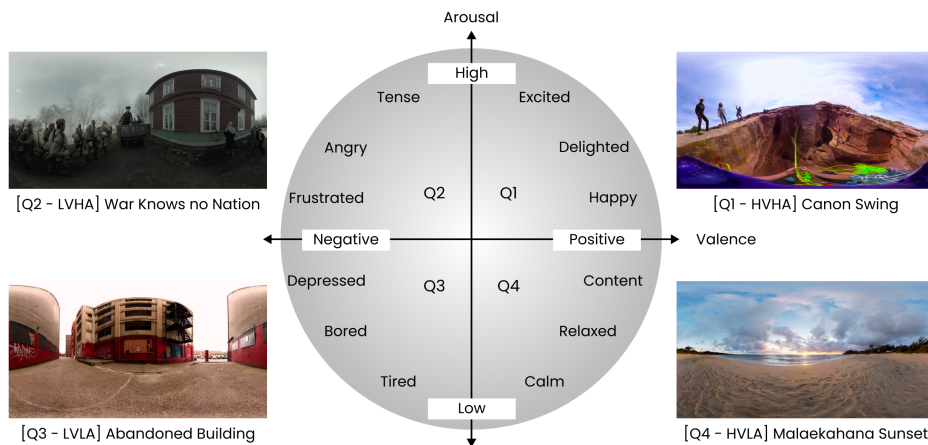


Figure 2: The figure presents still images extracted from 360° videos used in the experiment to display various environments to the participants. These images represent different valence-arousal combinations, including Q1: High-Valence-High-Arousal (HVHA), Q2: Low-Valence-High-Arousal (LVHA), and Q3: Low-Valence-Low-Arousal (LVLA), Q4: High-Valence-Low-Arousal (HVLA). The videos were selected from the publically available 360° VR video dataset (Li et al. (2017).)

## 29 1 EEVR Overview

30 The *EEVR* dataset comprises synchronized pairs of physiological signals and textual data. It includes  
 31 responses to four self-assessment questions regarding perceived arousal, valence, dominance, and  
 32 discrete emotions ratings collected using PANAS questionnaires (which were further utilized to  
 33 calculate Positive and Negative Affect Score). The *EEVR* dataset was collected using Virtual Reality  
 34 (VR) 360° videos as the elicitation medium. The videos utilized in the dataset were selected based on  
 35 their arousal and valence ratings to cover all four quadrants of the Russell circumplex emotion model  
 36 (Russell et al. (1989)), as shown in Figure 2. The remainder of the supplementary materials provide  
 37 detailed information about the *EEVR* dataset. Figure 3 provides a datasheet for the *EEVR* dataset  
 38 based on Gebru et al. (2018). The experiment setup is presented in Figure 1.

### 39 1.1 EEVR Size Details

40 The *EEVR* dataset consists of data from 37 healthy participants who agreed to share their data  
 41 publicly. Although 41 participants were involved in the data collection study, data from only  
 42 37 participants is publicly available. During data acquisition, the data from four participants was  
 43 damaged due to factors like motion sickness in the VR environment and issues with sensor attachment.  
 44 Consequently, the dataset provides physiological signal data (including Electrodermal Activity (EDA)  
 45 and Photoplethysmogram (PPG) signals) and textual descriptions of emotions felt during each  
 46 emotional stimulus for 37 participants.

47 The *EEVR* dataset comprises 296 tasks in total, with each participant experiencing eight VR 360°  
 48 videos shown to induce emotions from all four quadrants of the Russell emotion model (two videos  
 49 from each quadrant). Table 1 presents a summary of the minute durations for each video, along with  
 50 their respective playlist details. Further details regarding the playlist and video order are elaborated  
 51 in Section 4.2.

Video Name	Count (minutes)	Playlist
The Displaced	3:23	1, 3
Happyland 360	2:43	1, 3
Jailbreak 360	3:06	1, 3
War Knows No Nation	3:15	1, 3
Canyon Swing	1:44	1, 3
Redwoods Walk Among Giants	2:00	1, 3
Speed Flying	2:34	1, 3
Instant Caribbean Vacation	2:30	1, 3
The Nepal Earthquake Aftermath	3:09	2, 4
Zombie Apocalypse Horror	3:00	2, 4
Abandoned building	3:00	2, 4
Kidnapped	2:58	2, 4
Mega Coaster	1:57	2, 4
Malaekahana Sunrise	3:29	2, 4
Puppies host SourceFed for a day	1:20	2, 4
Great Ocean Road	1:58	2, 4

Table 1: Video names with their duration details and the playlist number

### 52 1.2 EEVR Organization and File formats

53 The *EEVR* dataset, as downloaded, is organized into two main subdirectories, as illustrated in Figure  
 54 4. The first subdirectory contains processed physiological data, including CSV files with raw EDA  
 55 and PPG data, organized by participant details and Video ID for ease of use. It also includes EDA  
 56 and PPG features CSV files extracted using the feature extraction pipeline. The second subdirectory  
 57 holds raw data and is divided into four playlist folders. Each playlist folder contains directories for  
 58 subject-wise raw EDA text files, raw PPG text files, annotation text files, and raw ACQ files in the

<h1>EEVR Dataset Facts</h1>	
<b>Dataset</b> EEVR	
Motivation	
<b>Summary</b> EEVR is a multimodal dataset designed for emotion recognition. It comprises physiological signal data collected from wearable sensors along with raw textual captions corresponding to each emotion elicitation segment.	
<b>Example Use Case</b> Emotion Recognition, Arousal classification, Valence classification, Personality Recognition	
<b>Original Authors</b> P. Singh, R. Budhiraja, A. Gupta, A. Goswami, M. Kumar, P. Singh	
MetaData	
<b>URL</b>	<a href="https://melangelabiiitd.github.io/EEVR/">https://melangelabiiitd.github.io/EEVR/</a>
<b>Keywords</b>	Emotion Recognition, Wearable sensor, Physiological signal
<b>Format</b>	.acq, .csv, .txt
<b>Ethical Review Approval</b>	IRB-IIIT-Delhi, NECRBHR
<b>License</b>	CC BY-NC-SA
<b>First Release Year</b>	2024
Sensors	
<b>EDA</b>	SS57LA, 4-channel Biopac MP36
<b>PPG</b>	SS4LA, 4-channel Biopac MP36
Data Annotation	
<b>Self-Assessments</b>	Arousal, Valence, Dominance, Positive Affect Score, Negative Affect Score, Familiarity, Discrete Emotions
<b>Textual Labels</b>	Raw Textual Description per video stimuli
<b>Additional Data</b>	Personality Score (BFI10), GHQ, VRSQ
Participants	
<b>Count</b>	37
<b>Gender</b>	21 males, 16 females
<b>Age</b>	18-33 (M=23.1, SD=4.02)
<b>Background</b>	Bachelor's (24), Master's (8), Senior High School (4), Doctorate (3)
Dataset Size	
<b>Total Size</b>	668 MB
<b>Physiological Data Duration</b>	797 minutes and 83 seconds

Figure 3: Dataset Summary card for EEVR, constructed based on (Geburu et al. (2018))

59 original Biopac format. All the physiological data files are organized in .CSV and .TXT formats,  
 60 making them easily usable for all programming languages. All physiological signals were initially  
 61 sampled at 2000Hz but were downsampled to 128Hz for PPG and 15.68Hz for EDA to reduce  
 62 computational requirements. The .ACQ files contain the original 2000Hz data, while the .TXT and  
 63 .CSV files are the downsampled versions used for experimentation. The downsampling frequencies  
 64 were chosen based on previous research to ensure no information was lost. Participant 29's data was  
 65 collected in two parts due to a disconnection during the experiment, resulting in two raw data files.

66 Additionally, there are three files: one detailing participant information collected during the data  
 67 collection process, second with self-assessment details collected during data collection, and third  
 68 with text data from semi-structured interviews for each participant-video segment. We have also  
 69 included VR\_Application.zip file containing the VR environment simulation build file and video  
 70 resources. The Participant\_details is organized into sheets for Participant Details, GHQ-12, and  
 71 BFI-10. The Self\_Assessment file is further organized into sheets for the Pre-exposure Questionnaire,  
 72 Post-exposure Questionnaire, Affect-Personality Score, and VRSQ Scores. The text data file also



73 contains sheets for Text-Labels (Text description with information on participant ID and video ID  
 74 and labels), Video description, and Video ID and Video Name mapping information.

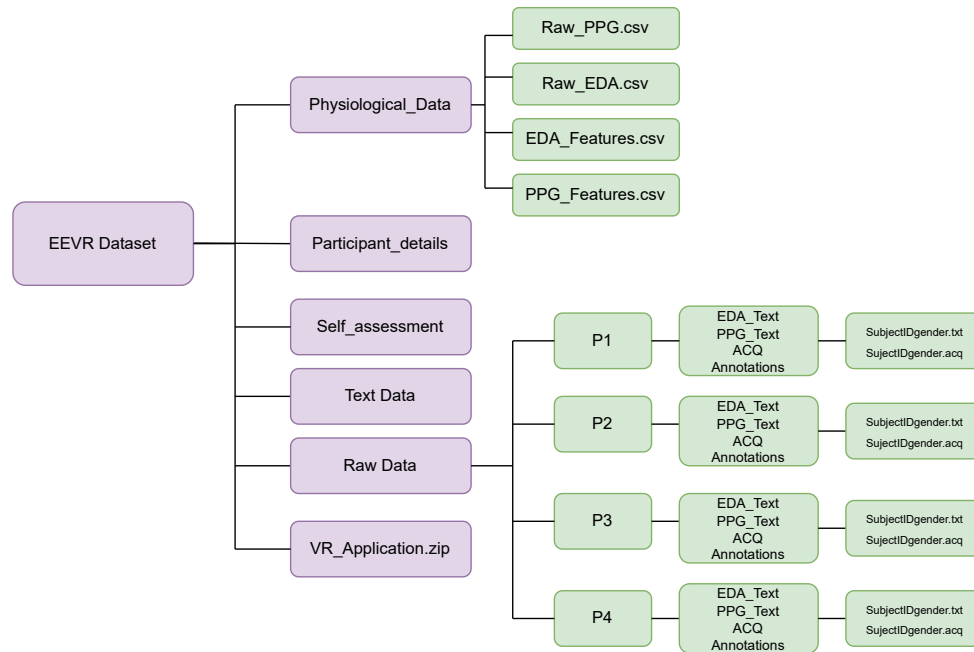


Figure 4: File Organization of EEVR dataset

## 75 2 EEVR Usage and Publishing

### 76 2.1 Intended Uses

77 The *EEVR* dataset is collected and published to further the research in Emotion recognition using  
 78 physiological signal data. The dataset is a resource for synchronised physiological signals, a textual  
 79 description of emotions felt and annotations in the form of perceived self-reports. Several use  
 80 cases, including extracting emotions from each modality and analyzing the correlations between  
 81 physiological signals and various emotion labels. Further, the dataset can be used for pre-training  
 82 physiological signal-based models using contrastive techniques for zero-shot classification of various  
 83 tasks like Valence classification, Arousal classification and stimuli-label-based emotion classification.

### 84 2.2 Ethical Consideration

85 We acknowledge that, despite all precautions, there is a possibility of the dataset being misused by  
 86 malicious users. The authors take full responsibility for any rights violations that may occur during  
 87 the data collection process or any related work. They are committed to taking necessary actions, such  
 88 as removing data that poses such issues, to address any problems that arise.

### 89 2.3 EEVR Licensing, Hosting, and Maintenance Plan

90 The *EEVR* dataset and its relevant code file are available for researchers to use further. The dataset  
 91 is available to use under CC BY-NC-SA license for non-commercial research. It can be accessed  
 92 by filling in the Dataset Access Request Form on our [website](#). Upon completing the form, users  
 93 can access the *EEVR* dataset stored on Google Drive. The authors will maintain the dataset files for  
 94 the long term, ensuring that the file structures remain unchanged. The *EEVR* website will also be

95 maintained for the long term, providing users with easy access to download the dataset. All the code  
96 files are available under the MIT open-source licence on [github](#).

### 97 **3 Human Subjects Considerations**

98 The EEVR dataset collection study has been approved by the Institution review board <sup>1</sup> of IIT-  
99 Delhi registered with the National Ethics Committee Registry for Biomedical and Health Research  
100 (NECRBHR). The participants for this study were recruited through email invitations. Before data  
101 collection began, all participants were introduced to the study protocol and its purpose. They were  
102 also informed about privacy concerns and any risks involved in the study. All the subjects participated  
103 on a voluntary basis. Additionally, all participants are made aware of our exclusion criteria. No  
104 participants with experience or a history of heart issues, heart arrhythmia, high blood pressure,  
105 medical conditions affecting equilibrium, visual or auditory impairments, neurological ailments,  
106 cognitive challenges, psychological issues, or diagnosed depression were recruited for this study.  
107 Further, participants with motion sickness issues were also excluded to avoid VR sickness discomfort  
108 on our subjects. The Ag/AgCl electrodes <sup>2</sup> used in our study have been proven in the past to adhere  
109 well to various types of skin surfaces. Further, we informed all participants to stop the experiments if  
110 they felt any discomfort. All the participant's data is pseudo-anonymized before being made publicly  
111 available.

## 112 **4 Data Collection Protocol**

### 113 **4.1 Experiment Instruction and Sensors Preparation**

114 Before starting the data collection, all participants were asked to sit comfortably in a chair. They were  
115 then informed about the study's purpose, which was to collect physiological data related to various  
116 emotions using VR 360° videos. Participants were instructed to report any discomfort or issues  
117 during the data collection and were informed to stop the experiment at any time in case of discomfort.  
118 The use of sensors and VR headsets was explained, and participants were given time to ask questions  
119 and express any concerns. Consent was then obtained from each participant. The preparation of  
120 wearable sensors involved attaching EDA (SS57LA) and PPG (SS4LA) sensor modules to the Biopac  
121 MP36 acquisition system. EL507 electrodes were prepared with isotonic gel and attached to the  
122 participants' index and middle fingers, while the PPG module was attached to the ring fingers. The  
123 EDA sensor module was calibrated by removing and reattaching one of the sensor heads. Following  
124 this, the sensors were checked for accurate readings. Upon confirmation of acquisition without any  
125 error, the experiment is started.

### 126 **4.2 Stimulus Selection and Playlist Preparation**

127 For this experiment, a total of 16 videos were selected from a publicly available 360° VR dataset  
128 containing 73 videos Li et al. (2017). To curate this subset, we applied a heuristic protocol based on  
129 the circumplex model of emotions Russell et al. (1989). We chose four videos from each category  
130 of the circumplex model, selecting those with the maximum distance from the origin to ensure a  
131 diverse and comprehensive representation. The 16 videos were then divided into two subgroups of  
132 eight videos each, as mentioned in Table 2. This decision was based on feedback from a pilot study  
133 (participants not included in the main study), the total experiment duration, and considerations to  
134 prevent participant fatigue or motion sickness from VR exposure. Alternate videos from each quadrant  
135 were paired to create two balanced video sets, ensuring normalization between the subgroups and a  
136 balanced experimental setting. After dividing the videos into subgroups, they were arranged in two  
137 different orders. In the first order, videos were organized based on their valence ratings, representing  
138 the degree of pleasantness or unpleasantness associated with an emotional state. The videos were

<sup>1</sup><https://irb.iiitd.edu.in/>

<sup>2</sup><https://www.biopac.com/product/eda-electrodes/>

139 arranged randomly in the second order, creating four playlists. The sorting technique employed  
 140 for the study involved arranging videos from low negative valence to high positive valence. The  
 141 four playlists are as follows- *Playlist1: VideoSet1 - Random Order*, *Playlist2: VideoSet1 - Valence*  
 142 *Sorted Order*, *Playlist3: VideoSet2 - Random Order*, and *Playlist4: VideoSet2 - Valence Sorted*  
 143 *Order*. The playlists are shown in Table 3. Participants were allocated playlists in a gender-balanced  
 144 manner through random assignment. The valence-sorted orders were designed to induce emotions  
 145 to transition smoothly between emotions, starting from positive emotions, then introducing more  
 146 intense emotions, and finally transitioning to negative emotions. The random order was inspired by  
 147 prior works that randomly showed videos.

CMA Quadrant	VideoSet1 [V,A]	VideoSet2 [V,A]
HVHA	Canyon Swing [5.38, 6.88], Speed Flying [6.75, 7.42]	Mega Coaster [6.17, 7.17], Puppies host SourceFed for a day [7.47, 5.35]
HVLA	Redwoods Walk Among Giants [5.79, 2.0], Malaekahana Sunrise [6.57, 1.57]	Instant Caribbean Vacation [7.2, 3.2], Great Ocean Road [7.77, 3.92]
LVHA	Jailbreak 360 [4.4, 6.7], Zombie Apocalypse Horror [3.2, 5.6]	War Knows No Nation [4.93, 6.07], Kidnapped [4.83, 5.25]
LVLA	The Displaced [2.18, 4.73], The Nepal Earthquake Aftermath [2.73, 3.8]	Happyland 360 [3.33, 3.4], Abandoned Building [4.39, 2.77]

Table 2: Videos categorized based on Valence (V), Arousal (A) rating

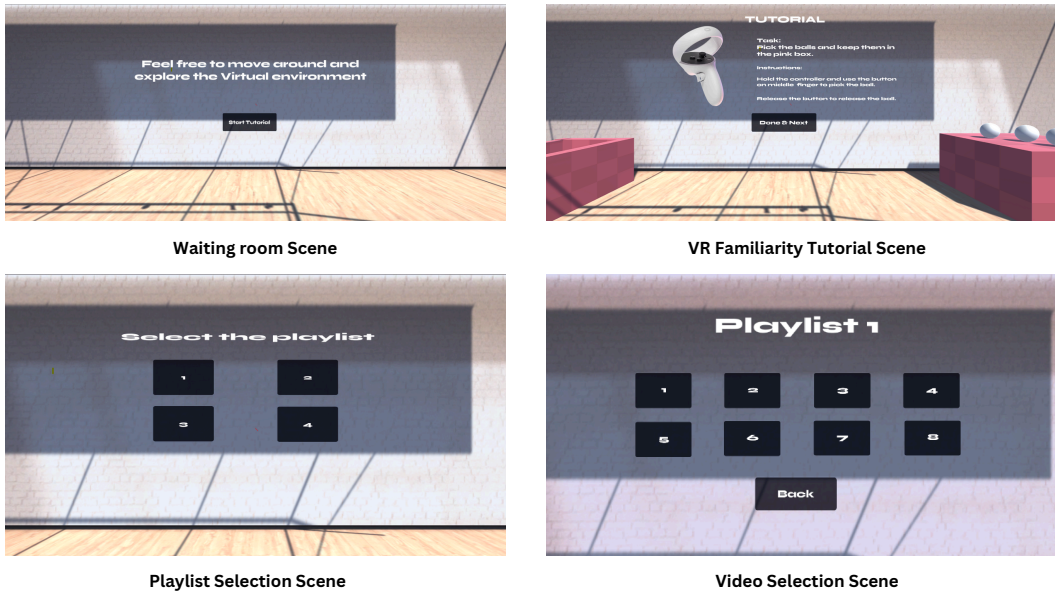


Figure 5: This figure illustrates the screens from Virtual Environment Room scene in following order: Waiting room scene, VR Familiarity Tutorial scene, Playlist Selection Scene and Video Selection Scene.

### 148 4.3 Virtual Reality Module Preparation

149 We developed a VR application for our experiment, enabling participants to experience emotionally  
 150 stimulating videos. The application consists of two main components: 1) an introductory module  
 151 to familiarize users with the VR environment and controllers, and 2) a video playback module for  
 152 presenting 360° videos. The application was created using Unity. Since many users were new to VR,  
 153 the introductory module starts with a "waiting room" scene designed to acclimate them to the VR

Video Name	Playlist ID-Video ID
The Displaced	P1V1
Happyland 360	P1V2
Jailbreak 360	P1V3
War Knows No Nation	P1V4
Canyon Swing	P1V5
Redwoods Walk Among Giants	P1V6
Speed Flying	P1V7
Instant Caribbean Vacation	P1V8
The Nepal Earthquake Aftermath	P2V1
Zombie Apocalypse Horror	P2V2
Abandoned Building	P2V3
Kidnapped	P2V4
Mega Coaster	P2V5
Malaekahana Sunrise	P2V6
Puppies host SourceFed for a day	P2V7
Great Ocean Road	P2V8
War Knows No Nation	P3V1
Redwoods Walk Among Giants	P3V2
Happyland 360	P3V3
Speed Flying	P3V4
Instant Caribbean Vacation	P3V5
Jailbreak 360	P3V6
The Displaced	P3V7
Canyon Swing	P3V8
Kidnapped	P4V1
Malaekahana Sunrise	P4V2
Zombie Apocalypse Horror	P4V3
Puppies host SourceFed for a day	P4V4
Great Ocean Road	P4V5
Abandoned Building	P4V6
The Nepal Earthquake Aftermath	P4V7
Mega Coaster	P4V8

Table 3: Video Names and Video ID

154 environment. Instructions, including text and images, are displayed on the walls using Unity’s XR UI  
155 canvas to guide users in interacting with and manipulating objects using the VR controller. To practice  
156 these skills, users complete a simple task of placing a ball in a bucket within the introductory scene.  
157 The second component allows users to experience 360° videos in VR. We curated four playlists,  
158 each containing eight videos, which were downloaded from a database<sup>3</sup> using the youtube-dl<sup>4</sup> tool.  
159 These videos are in equirectangular panoramic format with a 3840 x 2160 pixels resolution. In Unity,  
160 separate scenes were created for each video, with texture renderers mapping the video frames to a  
161 skybox surrounding a central camera. A script tracks video playback in each scene, and once a video  
162 finishes, the user is returned to the playlist menu to select another video. This setup allowed us to  
163 collect users’ physiological data for each video.

#### 164 4.4 Self-Assessment

165 Each task in our study is annotated using Valence, Arousal, and Dominance. Additional data on  
166 liking and familiarity was also collected using scales. The valence, arousal, dominance, and liking  
167 scales are presented in Figure 6. The familiarity was collected on a Likert scale of 1-5, with 1 being

<sup>3</sup><https://stanfordvr.com/360-video-database/>

<sup>4</sup><https://github.com/ytdl-org/youtube-dl>

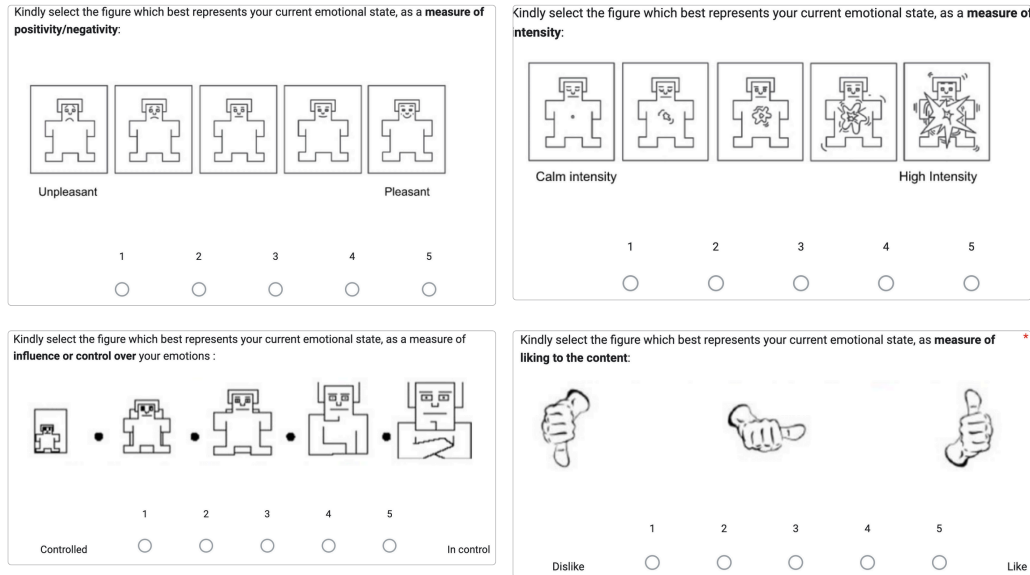


Figure 6: Illustration of self-assessment scales as following: Valence SAM, Arousal SAM, Dominance SAM, and Liking scale.

168 "very unfamiliar" and 5 being "very familiar". Similarly, PANAS scale annotations for ten positive  
 169 (Interested, Strong, Enthusiastic, Proud, Inspired, Determined, Alert, Attentive, Active) and ten  
 170 negative (Distressed, Irritable, Guilty, Scared, Upset, Hostile, Jittery, Ashamed, Nervous, Afraid)  
 171 emotions were also collected on a scale of 1-5, with 1 denoting "very slightly or not at all" to 5  
 172 denoting "extremely" for each emotion in the scale. To calculate the Positive Affect Score, we summed  
 173 up the scores for positive items (Interested, Strong, Enthusiastic, Proud, Inspired, Determined, Alert,  
 174 Attentive, and Active). This score can range from 10 to 50, with higher scores indicating higher  
 175 levels of positive affect. For the Negative Affect Score, we have added the scores for negative items  
 176 (Distressed, Irritable, Guilty, Scared, Upset, Hostile, Jittery, Ashamed, Nervous, Afraid). This score  
 177 also ranges from 10 to 50, with lower scores indicating lower levels of negative affect.

#### 178 4.4.1 GHQ-12

179 This study used a twelve-item General Health Questionnaire designed to measure non-psychotic  
 180 mental health. This scale is rated on a 4-point scale with a timeframe of "in the last one week." We  
 181 applied the Likert scoring method (0-1-2-3), where each of the four response options ("Not at all,"  
 182 "No more than usual," "Rather more than usual," "Much more than usual" or "Better than usual,"  
 183 "Same as usual," "Less than usual," "Much less than usual") is assigned a numerical value of 0, 1,  
 184 2, or 3. For each of the 12 questions, we summed the scores based on the responses given by the  
 185 respondents. The total score can range from 0 to 36. A lower total score (closer to 0) indicates better  
 186 mental health and lower psychological distress, while a higher total score (closer to 36) suggests  
 187 higher levels of psychological distress and potential mental health issues.

#### 188 4.4.2 Personality

189 The personality questionnaire, depicted in Figure 7 with item numbers, employs the BFI-10 (Big Five  
 190 Inventory-10) scoring method to derive personality scores. This method involves assigning scores  
 191 to each item based on the respondent's selection. For Extraversion, item 1 is reverse-scored (For  
 192 reverse-scoring item, subtract the respondent's original score from the highest possible score on the  
 193 scale plus one.), while item 5 is scored as is. In Agreeableness, item 2 is scored as is, and item 7 is  
 194 reverse-scored. Conscientiousness is determined by reversing the score for item 3 and scoring item  
 195 8 as is. Neuroticism involves reversing the score for item 4 and scoring item 9 as is. Openness to

I see myself as someone who ...	Disagree strongly	Disagree a little	Neither agree nor disagree	Agree a little	Agree strongly
1. ... is reserved	(1)	(2)	(3)	(4)	(5)
2. ... is generally trusting	(1)	(2)	(3)	(4)	(5)
3. ... tends to be lazy	(1)	(2)	(3)	(4)	(5)
4. ... is relaxed, handles stress well	(1)	(2)	(3)	(4)	(5)
5. ... has few artistic interests	(1)	(2)	(3)	(4)	(5)
6. ... is outgoing, sociable	(1)	(2)	(3)	(4)	(5)
7. ... tends to find fault with others	(1)	(2)	(3)	(4)	(5)
8. ... does a thorough job	(1)	(2)	(3)	(4)	(5)
9. ... gets nervous easily	(1)	(2)	(3)	(4)	(5)
10. ... has an active imagination	(1)	(2)	(3)	(4)	(5)

Figure 7: Illustration of BFI-10 personality scale used for our experiment with item number.

196 Experience is evaluated by reversing the score for item 5 and scoring item 10 as is. By applying these  
 197 scoring guidelines to each item, we calculate the total score for each trait.

#### 198 4.4.3 VRSQ

199 The VRSQ questionnaire is illustrated in Figure 8 with question numbers. To determine the VRSQ  
 200 score, we first calculated two sub-scores: A and B. Sub-score A is obtained by summing the  
 201 responses to questions 1 through 4, while sub-score B is derived from questions 5 through 9. Then, to  
 202 standardize these scores, A is divided by 12 and multiplied by 100 to yield C, and B is divided by 15  
 203 and multiplied by 100 to produce D. Finally, the VRSQ score is calculated as the average of C and D,  
 204 providing a comprehensive measure of VR sickness for an individual Kim et al. (2018).

1. General discomfort	<u>None</u>	<u>Slight</u>	<u>Moderate</u>	<u>Severe</u>
2. Fatigue	<u>None</u>	<u>Slight</u>	<u>Moderate</u>	<u>Severe</u>
3. Headache	<u>None</u>	<u>Slight</u>	<u>Moderate</u>	<u>Severe</u>
4. Eye strain	<u>None</u>	<u>Slight</u>	<u>Moderate</u>	<u>Severe</u>
5. Difficulty focusing	<u>None</u>	<u>Slight</u>	<u>Moderate</u>	<u>Severe</u>
6. Fullness of the Head	<u>None</u>	<u>Slight</u>	<u>Moderate</u>	<u>Severe</u>
7. Blurred vision	<u>None</u>	<u>Slight</u>	<u>Moderate</u>	<u>Severe</u>
8. Dizziness with eyes closed	<u>None</u>	<u>Slight</u>	<u>Moderate</u>	<u>Severe</u>
9. *Vertigo	<u>None</u>	<u>Slight</u>	<u>Moderate</u>	<u>Severe</u>

Figure 8: Illustration of Virtual Reality Sickness scale with questions as used in our experiments.

## 205 5 Data Analysis and Experiments

### 206 5.1 Content Analysis

207 In Figure 9, we illustrate the frequency distribution of self-reported annotations for each scale. Our  
208 analysis showed that the self-reports are mostly unbalanced. For example, the valence label tends  
209 to skew towards positive values, while the arousal label is predominantly neutral. Additionally,  
210 most participants reported a high level of control over their emotions on the dominance scale. Most  
211 participants also indicated that they liked the content used to induce emotions, which may explain the  
212 low negative affect scores across the board. Familiarity with the content was mostly high among the  
213 participants. Positive affect scores were more evenly distributed than skewed negative affect scores.  
214 Additionally, we found the GHQ scores of participants are evenly distributed. We observed that due  
215 to the subjective nature of emotions, participants' high levels of liking and perceived control over  
216 their emotions likely contributed to their overall positive reports.

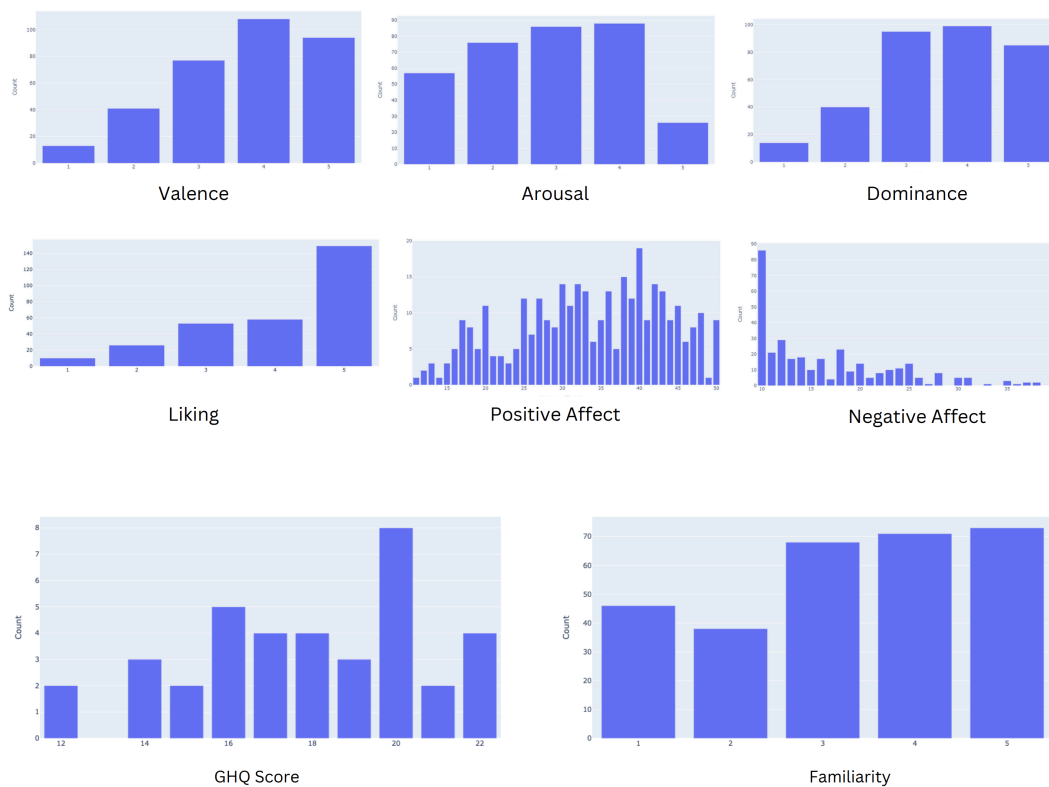


Figure 9: Frequency Distribution of self-reported annotations for Valence, Arousal, Dominance, Liking, Positive Affect, Negative Affect, GHQ Scores and Familiarity.

### 217 5.2 Data Cleaning

218 The physiological signal data was initially collected as ACQ files from Biopac, which allows the  
219 extraction of the data as text files. Before extraction, the signals were manually checked for errors,  
220 with any erroneous sections labeled for post-processing. The data was then downsampled using the  
221 software: EDA data to 15.625 Hz and PPG data to 125 Hz. The Biopac system was also used to  
222 calculate the BPM for the PPG data. After downsampling and BPM calculation, each participant's  
223 physiological signal text files and annotation files were downloaded. These files were then uploaded  
224 to Python using the pandas library and cleaned to remove all segments labeled as errors. The data was  
225 checked for NaN values and outliers, specifically PPG values outside the normal 35-140 BPM range  
226 and EDA values outside the 0-60  $\mu$ S range. Following this filtering, the signal data was labeled with



227 video ID, video name, playlist ID, and gender details based on timestamps. This helps us prepare the  
228 raw CSV files for further analysis.

### 229 5.3 Text Data Preparation

230 The textual descriptions were collected using a semi-structured interview technique, where an  
231 interviewer asked participants to explain their experiences qualitatively. Audio recordings were made  
232 for both the interviewer and interviewee. These recordings were then converted into text format using  
233 the Google Cloud Speech-to-Text API<sup>5</sup>. After conversion, the text data was manually checked for  
234 errors. Finally, the text data of participants' responses was extracted and compiled into a CSV file for  
235 further analysis.

### 236 5.4 Text Data Analysis

237 To assess the quality of our textual descriptions, we conducted a correlation analysis between these  
238 descriptions and the participant-reported valence and arousal (V/A) ratings (see Figure ??). For this  
239 analysis, we first extracted text embeddings using DistilBERT and subsequently applied Principal  
240 Component Analysis (PCA) to reduce the dimensionality of these embeddings. The resulting principal  
241 components were then visualized using a heatmap to illustrate their relationship with the V/A labels.  
242 Our findings indicate a strong correlation between the first principal component (PC1) and the V/A  
243 ratings, suggesting that the textual data closely aligns with the self-reported labels.

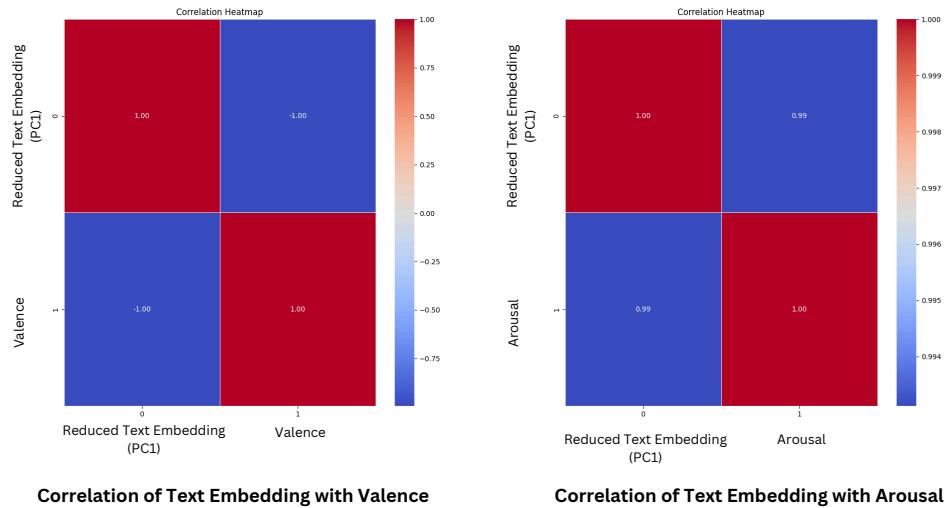


Figure 10: Illustration of correlation between V/A labels and textual descriptors

### 244 5.5 Physiological Features

245 For EDA data following signal cleaning and signal decomposition into tonic and phasic components,  
246 we have manually extracted the time domain features, such as statistical features, SCR-specific,  
247 and frequency domain features, such as power band features, variance, range, skewness, kurtosis.  
248 Similarly, we have extracted features using the Neurokit library for PPG data following the filtering  
249 and winsorization. We extracted Heart Rate (HR), Heart Rate Variability (HRV) Time-Domain  
250 Features, and Heart Rate Variability (HRV) Frequency-Domain Features. Following feature extraction,  
251 we analyzed correlation and dropped features with high correlation. Table 4 mentions the final features  
252 selected for classification.

<sup>5</sup><https://cloud.google.com/speech-to-text>



Signal	Selected Features
PPG	'BPM', 'IBI', 'PPG_Rate_Mean', 'HRV_MedianNN', 'HRV_Prc20NN', 'HRV_MinNN', 'HRV_HTI', 'HRV_TINN', 'HRV_LF', 'HRV_VHF', 'HRV_LFn', 'HRV_HFn', 'HRV_LnHF', 'HRV_SD1SD2', 'HRV_CVI', 'HRV_PSS', 'HRV_PAS', 'HRV_PI', 'HRV_C1d', 'HRV_C1a', 'HRV_DFA_alpha1', 'HRV_MFDFA_alpha1_Width', 'HRV_MFDFA_alpha1_Peak', 'HRV_MFDFA_alpha1_Mean', 'HRV_MFDFA_alpha1_Max', 'HRV_MFDFA_alpha1_Delta', 'HRV_MFDFA_alpha1_Asymmetry', 'HRV_ApEn', 'HRV_ShanEn', 'HRV_FuzzyEn', 'HRV_MSEn', 'HRV_CMSEn', 'HRV_RCMSEn', 'HRV_CD', 'HRV_HFD', 'HRV_KFD', 'HRV_LZC'
EDA	'ku_eda', 'sk_eda', 'dynrange', 'slope', 'variance', 'entropy', 'insc', 'fd_mean', 'max_scr', 'min_scr', 'nSCR', 'meanAmpSCR', 'maxAmpSCR', 'meanRespSCR', 'sumAmpSCR', 'sumRespSCR'

Table 4: List of Selected Features for Classification Tasks

## 253 5.6 Experiment Details and Results Analysis

### 254 5.6.1 Experimental Setup

255 We have used a machine with an AMD EPYC 7763 64-core Processor CPU and NVIDIA A100 40GB  
 256 GPU to train all our models. Training a classical machine learning model on physiological signals for  
 257 any classification task (Arousal, Valence, and Stimulus-Label) took around 10-15 minutes of CPU  
 258 time. Similarly, training the BERT-based text classification models took approximately 20 minutes of  
 259 GPU time for each classification task. The Contrastive Language-Signal Pre-training (CLSP) Model  
 260 required around 55.5 GPU hours for training 7 epochs in a Leave-One-Subject-Out (LOSO) setup for  
 261 37 participants. We trained the models for Electrodermal Activity (EDA) and Photoplethysmogram  
 262 (PPG) signals separately, totaling 333 GPU hours for training 7 epochs across all tasks (Stimulus  
 263 Label, Valence Label, and Arousal Label). Due to the CPU-intensive nature of our CLSP experiments  
 264 and the limited CPU computing power available, our experiments took longer than expected.

### 265 5.6.2 Discussion on Physiological Baseline

266 Our baseline results for the Valence and Stimulus\_Label classification task across all classical  
 267 machine-learning models were better than random, suggesting that our models are able to separate  
 268 features for these labels. We observed that stimulus labels are easy to predict using physiological  
 269 signal-based features as compared to subjective labels like valence and arousal. The PPG+EDA  
 270 features gave the best performance for valence classification. And EDA features provided the  
 271 best performance for Stimulus\_Label classification. The results were nearly random for Arousal  
 272 classification even after performing duplicate upsampling, suggesting arousal classification is a  
 273 difficult label to predict purely based on physiological signals-based features. We have visualized our  
 274 EDA and PPG handcrafted features for all three tasks using t-SNE algorithms as shown in Figure  
 275 5.6.2. We observed that t-SNE features are separable for Stimulus\_Label in the case of EDA data.  
 276 While for other labels the features are overlapping. This suggests the need for more complex models  
 277 and better representation learning for valence and arousal prediction.

### 278 5.6.3 Discussion on Text Baseline

279 To assess the quality of our text data, we performed all three classification tasks using only text data  
 280 as our input. Text-based classification models significantly outperformed classification models trained  
 281 on only physiological signals. This better performance is likely due to the use of large pre-trained  
 282 embedding models. We used the DistilBERT and XLM-RoBERTa Base for classification, where  
 283 DistilBERT performed better. We trained all models with a batch size of 16 over 7 epochs and a  
 284 learning rate of  $2e-5$ . Furthermore, we visualized the embeddings from our models and found that  
 285 they are visually distinct, as illustrated in Figure 12.

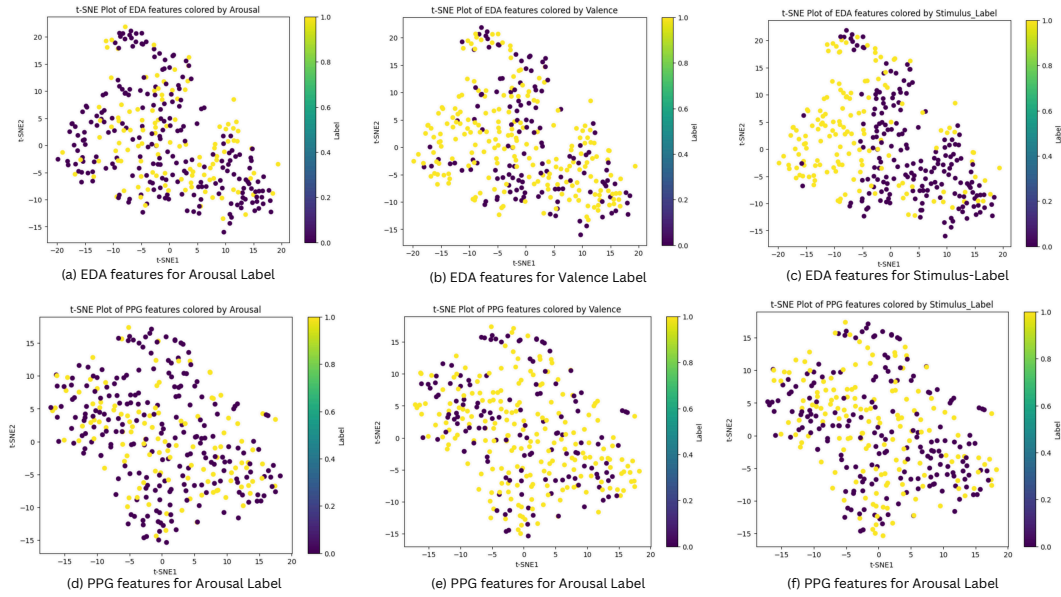


Figure 11: t-SNE plot depicting feature distribution of physiological signals according to various labels (Arousal, Valence, and Stimulus-Label).

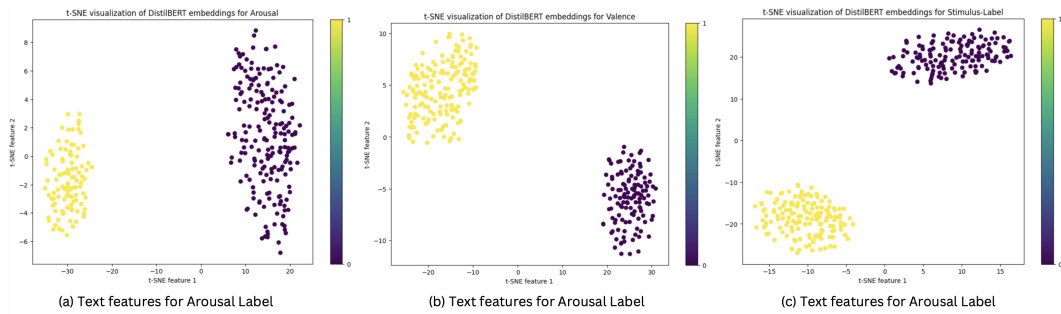


Figure 12: t-SNE plot depiction of Text data features for our three labels: Arousal, Valence, and Stimulus-Label

286 **5.6.4 Discussion on CLSP**

287 To assess the significance of aligning textual descriptions to physiological signal data, we performed  
 288 CLSP training for all three labels using EDA only, PPG only, and EDA+PPG handcrafted features  
 289 (for 296 tasks excluding baseline) along with text data. To compare our CLSP results, we first  
 290 trained a hand-crafted features-based neural network (HC+NN) model with two hidden linear layers  
 291 of dimensions 50 and 100. These models were trained for a batch size of 32 with 200 epochs,  
 292 using a learning rate of 0.001. For optimization, we utilized the Adam optimizer with beta values  
 293 of 0.9 for beta1 and 0.999 for beta2 and an epsilon value of 1e-8. The CLSP model was then  
 294 trained for contrastive objectives using the HC+NN model for physiological signal embeddings  
 295 and the DistillBert model with a projection head of a single linear layer with dimension 100 for  
 296 text embeddings. The project head was added to match the dimensionality of the text embeddings  
 297 (originally 768) with signal embedding. The training was conducted for 15 epochs with a learning  
 298 rate of 0.001 and batch size of 32. We found that our results were significantly better for arousal  
 299 and valence labels, suggesting the importance of augmenting text data for training subjective labels  
 300 like Arousal and Valence. The results were not as good for Stimulus\_Labels, indicating that the  
 301 features extracted from signals do not complement the features extracted from text. This mismatch

302 could have happened due to the subjective nature of textual descriptions that Stimulus\_Labels cannot  
303 capture. This misalignment might have led to ineffective contrastive training. The EDA+Text  
304 outperformed PPG+Text and PPG+EDA+Text, suggesting that EDA features might align more with  
305 textual descriptions for arousal and valence labels. We also observed that the early fusion of EDA and  
306 PPG features for CLSP has led to poor performance compared to EDA-only and PPG-only features,  
307 with Text indicating a need for designing a more complex fusion technique. Overall, our results  
308 suggest that aligning text data with physiological signals can improve learning for subjective labels,  
309 achieving results that cannot be attained with objective labels alone.

## 310 **References**

- 311 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M. Wallach,  
312 Hal Daumé III, and Kate Crawford. 2018. Datasheets for Datasets. *CoRR* abs/1803.09010 (2018).  
313 arXiv:1803.09010 <http://arxiv.org/abs/1803.09010>
- 314 Hyun K Kim, Jaehyun Park, Yeongcheol Choi, and Mungyeong Choe. 2018. Virtual reality sickness  
315 questionnaire (VRSQ): Motion sickness measurement index in a virtual reality environment.  
316 *Applied ergonomics* 69 (2018), 66–73.
- 317 Benjamin J Li, Jeremy N Bailenson, Adam Pines, Walter J Greenleaf, and Leanne M Williams. 2017.  
318 A public database of immersive VR videos with corresponding ratings of arousal, valence, and  
319 correlations between head movements and self report measures. *Frontiers in psychology* 8 (2017),  
320 2116.
- 321 James Russell, Anna Weiss, and G. Mendelsohn. 1989. Affect Grid: A Single-Item Scale of  
322 Pleasure and Arousal. *Journal of Personality and Social Psychology* 57 (Sept. 1989), 493–502.  
323 <https://doi.org/10.1037/0022-3514.57.3.493>